# Natural Language Processing Project 3

Joel, Janet, Clement
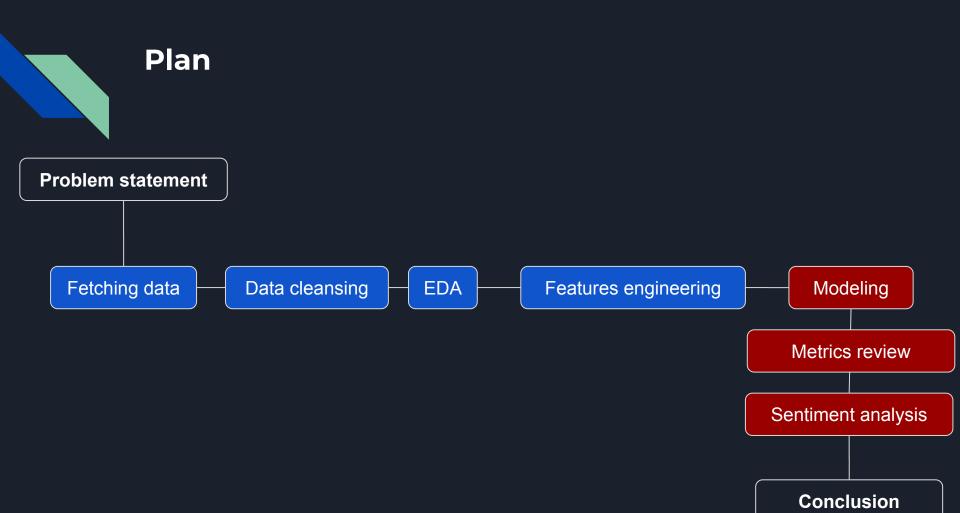
# Problem statement
*Can reddit to become an aggregator of conversations*

Reddit now want to display organic conversations from its members, but also any publicly available conversation on the internet, and classified it under one of its existing subreddits.

They have hired us to run a "stage 0 testing" to evaluate whether such an algorithm would even be effective based sources from reddit itself, let alone importing from other sources.

# Plan

**Problem statement**

Fetching data — Data cleansing — EDA — Features engineering — Modeling

Metrics review

Sentiment analysis

**Conclusion**

# Getting the data

Using Reddit API

**Size**: each subreddit got scrapped 200 post x 30 iterations = **6000 posts**

**Data pre-EDA cleansing**

- **Transform to lower**
- **Remove** empty posts
- **Handle** emojis

# EDA

1. **Most common words** in each df
   - Tokenizing
   - Lemmatizing & Stemming   (9498 columns vs 7381 columns)
   - Removing  stop words
   - Checking similarity: Unigrams vs Bi-grams

2**. Counting words**
   - Check for min, max and mean word count on the individual datasets
   - Dropping the rows that contains no words in titles  (Users include an emoji or punctuation in place of a word in title)
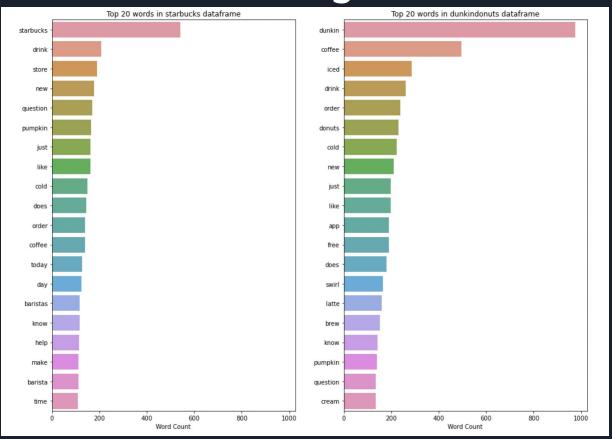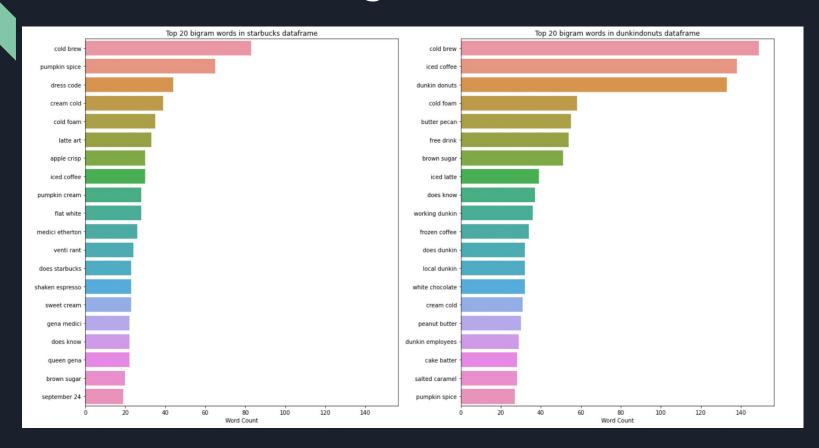   - Plotting in WordClouds

3. **Posts length** in each df
Is there a pattern in length that may influence the model?

| | title | subreddit | words_in_sentence |
|---|---|---|---|
| **327** | . | starbucks | 0 |
| **506** | 😌 | starbucks | 0 |
| **1218** | 🙂 | starbucks | 0 |
| **1333** | . | starbucks | 0 |
| **2318** | 🔲 | starbucks | 0 |

# Visualization on unigram



Top 20 words in starbucks dataframe

| Word | Approx. Word Count |
|------|-------------------|
| starbucks | ~550 |
| drink | ~210 |
| store | ~190 |
| new | ~180 |
| question | ~170 |
| pumpkin | ~165 |
| just | ~160 |
| like | ~160 |
| cold | ~150 |
| does | ~145 |
| order | ~140 |
| coffee | ~140 |
| today | ~125 |
| day | ~120 |
| baristas | ~115 |
| know | ~115 |
| help | ~110 |
| make | ~110 |
| barista | ~105 |
| time | ~100 |

Top 20 words in dunkindonuts dataframe

| Word | Approx. Word Count |
|------|-------------------|
| dunkin | ~970 |
| coffee | ~490 |
| iced | ~290 |
| drink | ~260 |
| order | ~240 |
| donuts | ~230 |
| cold | ~220 |
| new | ~210 |
| just | ~200 |
| like | ~200 |
| app | ~190 |
| free | ~190 |
| does | ~180 |
| swirl | ~165 |
| latte | ~160 |
| brew | ~150 |
| know | ~140 |
| pumpkin | ~140 |
| question | ~135 |
| cream | ~130 |

# Visualization on Bigrams



Top 20 bigram words in starbucks dataframe

Top 20 bigram words in dunkindonuts dataframe
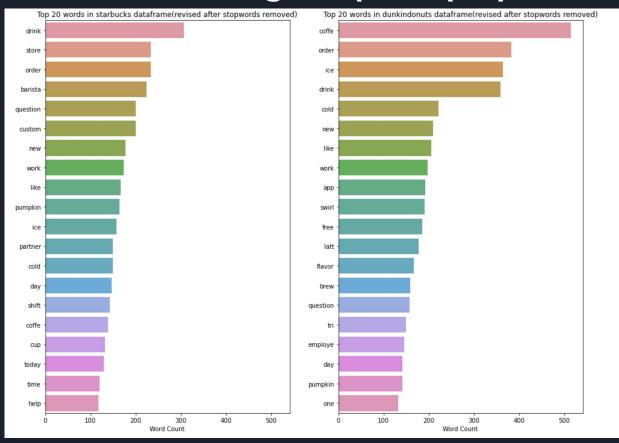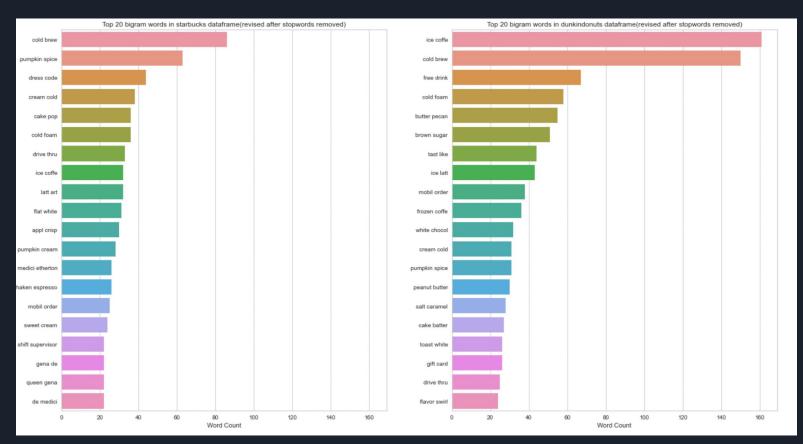
# WordCloud on raw title text

# Preprocessing

- Converting text in title to lowercase
- Removing numbers
- Removing URLs
- Removing punctuation
- Applying stemming
- Remove stop words, including brand names and other miscellaneous words like 'tall', 'grande', 'venti' all of which describes the cup size of Starbucks
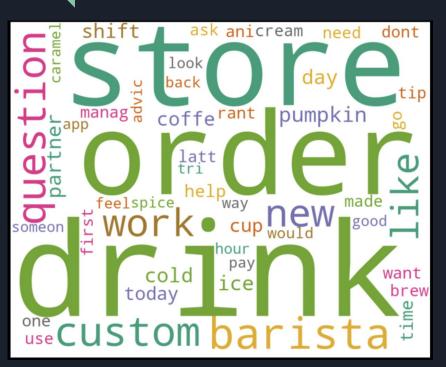
# Visualization on unigram (After preprocessing)



Top 20 words in starbucks dataframe(revised after stopwords removed)

Top 20 words in dunkindonuts dataframe(revised after stopwords removed)

# Visualization on Bigrams (After preprocessing)



Top 20 bigram words in starbucks dataframe(revised after stopwords removed)

Top 20 bigram words in dunkindonuts dataframe(revised after stopwords removed)

# WordClouds post-stop words removal

# Classification model metrics

- Choices : Accuracy, Recall, Precision, F1

- Chosen: Accuracy
  - Target variable is balanced (0.50 each)
  - Consequences of False Positive and False Negative is more or less the same
  - No prioritizing of True Positive/ Negative or  False Positive/ Negative required in this case
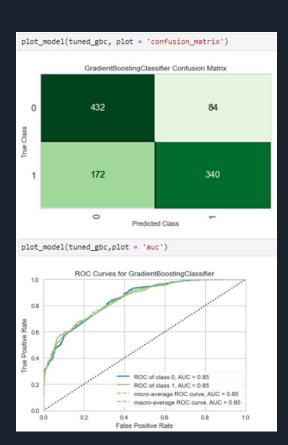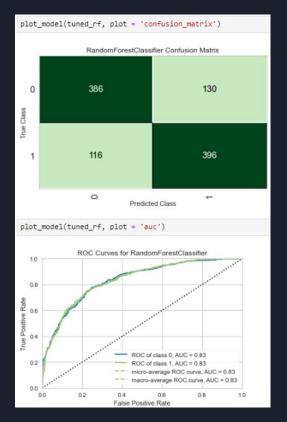
# Modeling

- Baseline model = Null model = 0.50

- Compare all models in PyCaret Library

- Top models:
  - Logistic Regression
  - Gradient Boosting Classifier
  - Extra Trees Classifier

- Tune the hyperparameter

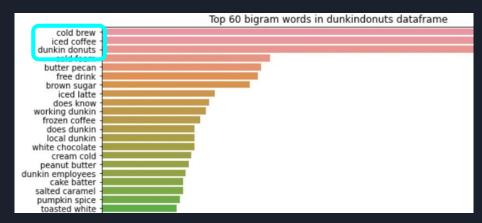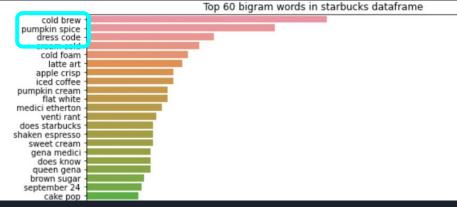| | Model | Accuracy | AUC | Recall | Prec. | F1 |
|---|---|---|---|---|---|---|
| lr | Logistic Regression | 0.7722 | 0.8535 | 0.7553 | 0.7808 | 0.7672 |
| gbc | Gradient Boosting Classifier | 0.7580 | 0.8542 | 0.6851 | 0.8025 | 0.7379 |
| et | Extra Trees Classifier | 0.7576 | 0.8314 | 0.7487 | 0.7619 | 0.7548 |
| ridge | Ridge Classifier | 0.7547 | 0.0000 | 0.7479 | 0.7571 | 0.7519 |
| rf | Random Forest Classifier | 0.7447 | 0.8351 | 0.7521 | 0.7406 | 0.7457 |
| svm | SVM - Linear Kernel | 0.7376 | 0.0000 | 0.7227 | 0.7442 | 0.7325 |

# Models EDA

# Final model and Evaluation

- Logistics Regression gives best accuracy of 0.78.
  - Well-fitted
  - Interpretable
  - Computational cheaper
  - Suitable for binary classification problem

Score Summary

| Model | Logistic Regression | Random Forest Classifier | Support Vector Machine | Gradient Boosting Classifier |
|-------|---------------------|--------------------------|------------------------|------------------------------|
| CV score | 0.7722 | 0.7447 | 0.7378 | 0.7580 |
| Train score | 0.7763 | 0.7607 | 0.7529 | 0.7510 |
| Test score | 0.7808 | 0.7726 | 0.7502 | 0.7611 |

# Sentiment Analysis - on Hot Topics

- Dunkin':
  - Cold brew
  - Iced coffee
  - Dunkin donuts
- Starbucks:
  - Cold brew
  - Pumpkin spice
  - Dress code
- Both:
  - Reward
  - Service


Top 60 bigram words in dunkindonuts dataframe
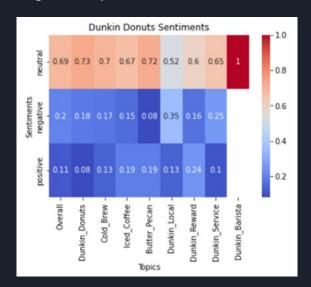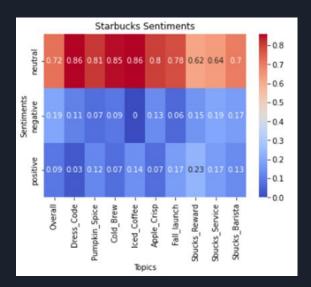

Top 60 bigram words in starbucks dataframe

# Sentiments

Majority posts seem neutral - Questions, Discussions

Well-received topics : Rewards (free drinks, free beverages, points)

Negative topics: Services (mobile order, app, staffs)

# Project conclusion

We have obtained a **reliable classifier to differentiate between 2 posts on reddit;**

**Recommendation:**
**We can use Logistic Regression as our model to allocate subreddits with satisfying accuracy.**

**We can add on Sentiment Analysis feature** to ease the navigation inside a subreddit for **HR | marketing business development teams**

**However**
To accurately classify a post coming from anywhere on the internet **will prove way more difficult**

# Future works

**Gathering more data** - fine tune model and enable time-series analysis

**Text-pre-processing steps**

**Improving on our models training** ( Gridsearch/ vectorizers etc...)

**Stacking models** to make predictions that have better performance than any single model in the ensemble